

Week 4 Lab - Sentiment Analysis II

Clarissa Boyajian

2022-04-24

```
raw_tweets <- read.csv("https://raw.githubusercontent.com/MaRo406/EDS_231-text-sentiment/main/dat/IPCC_
                      header = TRUE)

# load sentiment lexicons
bing_sent <- get_sentiments('bing')
nrc_sent <- get_sentiments('nrc')
```

Question 1: Clean Data

```
clean_data <-
  raw_tweets[, c(4, 6, 10:11)]

tweets <-
  tibble(id = seq(1:length(clean_data$Title)),
         text = clean_data$Title,
         date = as.Date(clean_data$Date, '%m/%d/%y'),
         sentiment = clean_data$Sentiment,
         emotion = clean_data$Emotion) %>%
  mutate(text = str_replace(string = text,
                             pattern = "http.*[:space:]",
                             replacement = ""),
         text = str_replace(string = text,
                             pattern = "http.*$",
                             replacement = ""),
         text = str_replace(string = text,
                             pattern = "@.*[:space:]",
                             replacement = ""),
         text = str_replace(string = text,
                             pattern = "@.*$",
                             replacement = ""),
         text = str_to_lower(text))

words <- tweets %>%
  select(id, date, text) %>%
  unnest_tokens(output = word,
                input = text,
                token = "words") %>%
  anti_join(stop_words, by = "word") %>%
  left_join(bing_sent, by = "word") %>%
  left_join(
```

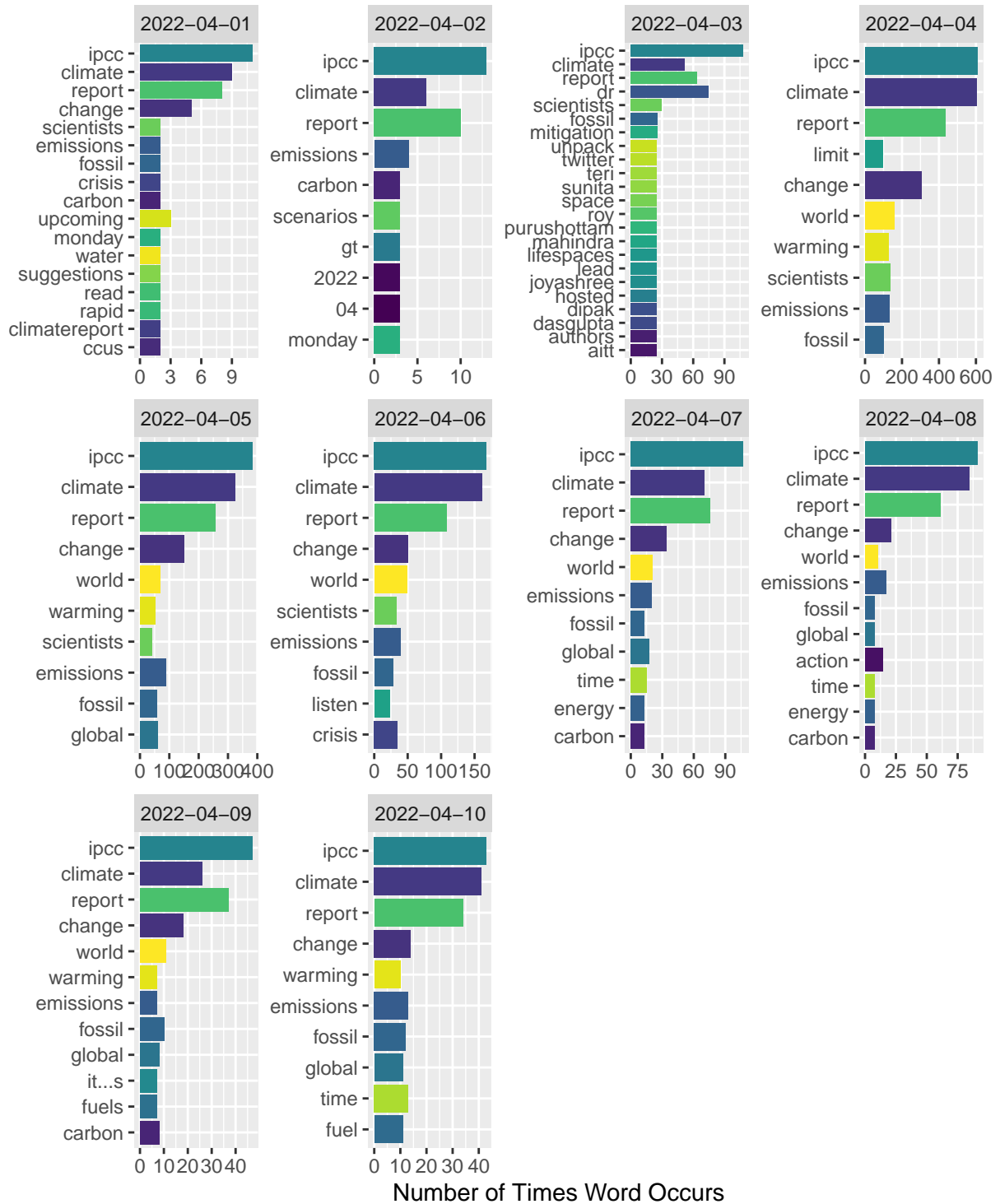
```
tribble(
  ~sentiment, ~sent_score,
  "positive", 1,
  "negative", -1),
by = "sentiment")
```

Question 2: Compare 10 Most Common Words per Day

```
words_common <- words %>%
  group_by(date, word) %>%
  summarise(count = n()) %>%
  group_by(date) %>%
  slice_max(count, n = 10)

ggplot(data = words_common,
  aes(x = count,
      y = reorder(word, count))) +
  geom_col(aes(fill = word)) +
  facet_wrap(~date, scales = "free") +
  guides(fill = "none") +
  scale_fill_viridis_d() +
  labs(x = "Number of Times Word Occurs",
      y = "",
      title = "Top 10 Words per Day")
```

Top 10 Words per Day



Question 3: Add Color to Wordcloud

```
words %>%
  count(word, sentiment, sort = TRUE) %>%
  acast(word ~ sentiment, value.var = "n", fill = 0) %>%
  comparison.cloud(colors = c("darkgreen", "red", "lightgrey"),
    max.words = 100)
```



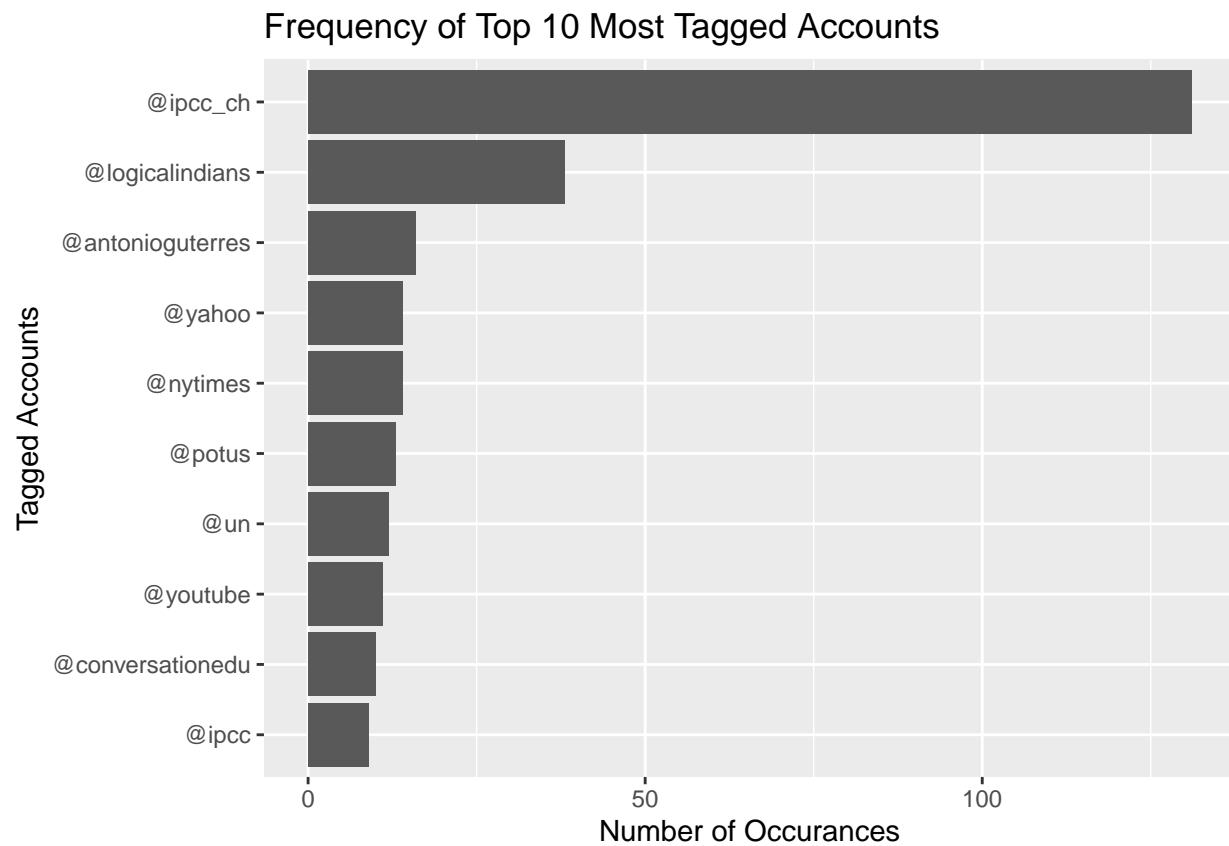
Question 4: Most Tagged Accounts

```
corpus <- corpus(clean_data$Title)

tagged_tweets <- tokens(corpus, remove_punct = TRUE) %>%
  tokens_keep(pattern = "@*") %>%
  dfm() %>%
  textstat_frequency(n = 10)

tagged_tweets_tidy <- tokens(corpus, remove_punct = TRUE) %>%
  tokens_keep(pattern = "@*") %>%
  dfm() %>%
  tidy()
```

```
ggplot(data = tagged_tweets,
       aes(x = frequency,
           y = reorder(feature, frequency))) +
geom_col() +
labs(x = "Number of Occurances",
     y = "Tagged Accounts",
     title = "Frequency of Top 10 Most Tagged Accounts")
```



Questeion 5